

539,493

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property
Organization
International Bureau



17 JUN 2005



(43) International Publication Date
1 July 2004 (01.07.2004)

PCT

(10) International Publication Number
WO 2004/055688 A1

(51) International Patent Classification⁷: G06F 15/80

(21) International Application Number:
PCT/GB2003/005532

(22) International Filing Date:
17 December 2003 (17.12.2003)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
0229368.6 17 December 2002 (17.12.2002) GB

(71) Applicant (for all designated States except US): ASPEX
TECHNOLOGY LIMITED [GB/GB]; Denmark House,
Denmark Street, High Wycombe, Buckinghamshire HP11
2ER (GB).

(72) Inventors; and

(75) Inventors/Applicants (for US only): JALOWIECKI,
Ian [GB/GB]; Denmark House, Denmark Street, High
Wycombe, Buckinghamshire HP11 2ER (GB). WHIT-
TAKER, Martin [GB/GB]; Denmark House, Denmark
Street, High Wycombe, Buckinghamshire HP11 2ER

(GB). LANCASTER, John [GB/GB]; Denmark House,
Denmark Street, High Wycombe, Buckinghamshire HP11
2ER (GB). BOUGHTON, Donald [GB/GB]; Denmark
House, Denmark Street, High Wycombe, Buckinghamshire
HP11 2ER (GB).

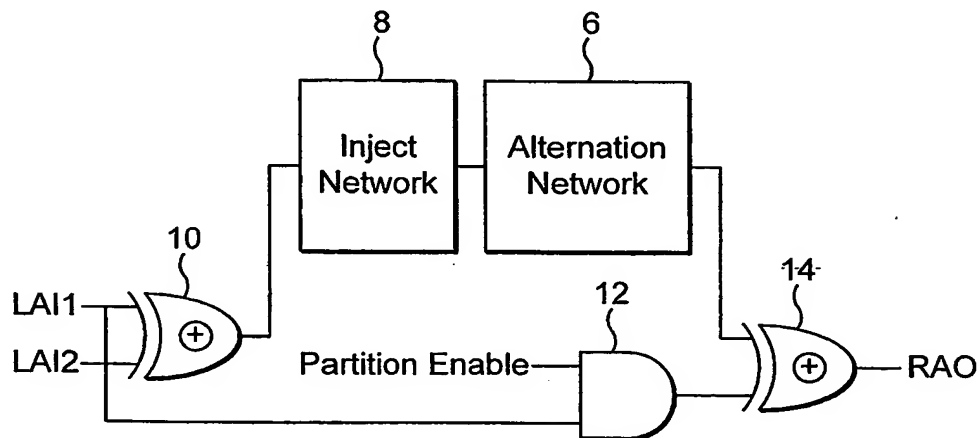
(74) Agents: AHMAD, Shelkh, Shakeel et al.; David Keltie
Associates, Fleet Place House, 2 Fleet Place, London
EC4M 7ET (GB).

(81) Designated States (national): AE, AG, AL, AM, AT, AU,
AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR,
CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD,
GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR,
KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN,
MW, MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU,
SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA,
UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

(84) Designated States (regional): ARIPO patent (BW, GH,
GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW),
Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM),
European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE,
ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE,
SI, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA,
GN, GQ, GW, ML, MR, NE, SN, TD, TG).

[Continued on next page]

(54) Title: SCALABLE PROCESSING NETWORK FOR SEARCHING AND ADDING IN A CONTENT ADDRESSABLE
MEMORY



(57) Abstract: An alternation network for use with a content addressable memory for implementing a divide and conquer algorithm is described. The alternation network comprises: a plurality of alternation modules connected in series together, each module comprising: a plurality of cascaded logic gates arranged to propagate a match parity signal via the gates along at least part of a matching result vector, the matching result vector being generated by execution of a matching instruction on the content addressable memory, and the logic gates being configured to change the parity of the match parity signal in accordance with the matching result vector; and a vector output arranged to output a parity level vector of the propagated match parity signal present at the each gate of the plurality of logic gates; a logic network for dividing the matching result vector into an odd match vector and an even match vector representing respectively odd and even numbered elements of the matching result vector, by use of the parity level vector; and means for writing a selected one of the odd and even match vectors to the content addressable memory.

WO 2004/055688 A1

WO 2004/055688 A1



Published:

— with international search report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

SCALABLE PROCESSING NETWORK FOR SEARCHING AND ADDING
IN A CONTENT ADDRESSABLE MEMORY

This invention relates to scalable processing networks, in particular but not restricted to processing networks employing concurrent microprocessor technology used in the communications industry.

A common type of processor function is to locate in a memory store all occurrences of a desired value, and then to determine the sum of all occurrences. Typically this is achieved through a series of synchronous matching actions wherein each memory store location is searched and matched in turn. However, this implementation is wasteful of time and processor cycles.

In a case where a processor has an associated local memory store, for example the multiprocessor implementation of Single Instruction Multiple Data (SIMD) processors which is the subject of the applicant's European Patent published as EP 0983556, the search and summation function as described in the above paragraph may be carried out through a type of inter-processor communication known as 'alternation'.

Efficient inter-processor communications in SIMD data parallel processors is paramount to performance. Algorithms that require efficient inter-processor communications are the so-called, binary 'divide-and-conquer' algorithms, whereby the multitude of processors in a data parallel architecture are subdivided into even/odd subsets in an efficient manner, enable the efficient processing of large amounts of data.

Figure 1 illustrates this principal, showing the summation of eight numbers (A to H) by successively adding even-odd pairs in $\log_2 N$ steps, where $N=8$. Thus the summation of A to H is achieved in three steps, compared to at least eight steps using the previous method of synchronous summation of each data value (A to H) in turn.

There have been some attempts in the past to implement these divide and conquer algorithms in hardware as alternation networks. These attempts have resulted in a basic building block being known which can be replicated to provide a scalable solution.

5

The basic logic building block of the known alternation network comprises a linear string of EXOR (Exclusive-OR) gates, one string per SIMD processor, connected so that the control signal is local to the processor, and the communication signal propagates to the next processor.

10

Figure 2 shows such a solution. The alternation network consists of EXOR gates 2 and a pair of AND gates 4a and 4b connected in the arrangement as shown. Upon instruction from a CPU (not shown in Figure 2) each SIMD processor generates a status signal called a 'match' signal, where the match signal from the n th processor is termed M_n . The application of the match signals causes the generation of a signal which propagates between processors, swapping state whenever it encounters a condition where a local match (M_n) is true. By combining a propagation term RAO in true or complement form with the original match signal, then two new match states can be generated from AND gates 4a and 4b, termed $M_{n_{\text{even}}}$ and $M_{n_{\text{odd}}}$. As can be seen from Figure 2, an inverter is provided at one input of each AND gate 4b such that only one output from the pair of AND gates 4a or 4b will be true at any given time with the same inputs going into both gates 4a and 4b. The propagation term RAO is then fed, for example, into the input of the next alternation network. The input LAI can be user defined or it can be an assumed value.

20

25

30

The versatility of this solution is that it allows the even and odd match sets to be easily generated, offering application flexibility. However, there is a significant drawback namely its very poor speed of execution, which is due to the propagating signal being limited by the gate delay of a single EXOR gate and the number of processors in the string.

Providing a prior art example, an alternation network comprising the VASP-256 data parallel processor comprises a network of 256 EXOR gates, one per processor. Execution of an alternate instruction using the prior art alternation network requires approximately eleven clock cycles (at 50MHz) to completely
5 generate an even-odd partition of M_0 to M_{255} . This as will be described later on is very slow in relative terms.

It is thus a first object of the invention to provide improved hardware support in microprocessor architecture for a binary alternation operation. Further objects
10 of the invention include provision of transparent operation across chip boundaries in modular multi-chip solutions with a minimum of external hardware support, and for a microprocessor possessing suitable architecture to complete the alternation process in a minimum number of clock cycles, preferably one cycle, irrespective of the number of microprocessors used.

15

Another object of the present invention is to implement a relatively fast asynchronous addition of remotely located elements in a content addressable memory.

20 It is therefore provided as a solution to the above problems the following an alternation network for use with a content addressable memory for implementing a divide and conquer algorithm, the network comprising: a plurality of alternation modules connected in series together, each module comprising: a plurality of cascaded logic gates arranged to propagate a match
25 parity signal via the gates along at least part of a matching result vector, the matching result vector being generated by execution of a matching instruction on the content addressable memory, and the logic gates being configured to change the parity of the match parity signal in accordance with the matching result vector; and a vector output arranged to output a parity level vector of
30 the propagated match parity signal present at the each gate of the plurality of logic gates; a logic network for dividing the matching result vector into an odd match vector and an even match vector representing respectively odd and even numbered elements of the matching result vector, by use of the parity

level vector; and means for writing a selected one of the odd and even match vectors to the content addressable memory.

The advantage of this structure is that it enables implementation of a divide and conquer algorithm to be realised in hardware such that massive improvements in speed can be achieved. The speedup is due to the ability to split the matching vector generated by the content addressable memory into odd and even vectors in a few clock cycles and thereafter to iteratively apply the search, and split process to arrive at a single processed element.

Further optional features to the invention in the above embodiments are described in the appended claims.

According to another aspect of the present invention there is provided a communications network for use with a content addressable memory for moving data between remote locations, the network comprising: dividing means for dividing a matching result vector generated by execution of a matching instruction on the content addressable memory, into an odd match vector and an even match vector representing respectively odd and even numbered elements of the matching result vector; assigning means for assigning the odd and even match vectors as source and destination locations; a plurality of cascaded logic gates arranged to propagate data via the gates along at least part of the matching result vector; input means for inputting a bit of a stored data value associated with a source location into the plurality of cascaded logic gates; applying means for applying the source and destination results to the plurality of cascaded logic gates to propagate the bit of stored data at the source location to the destination location; and extracting means for extracting the input bit from the plurality of cascaded logic gates at the destination location.

This communications network supports a search and add process which can be used effectively in a content addressable memory to carry out asynchronous remote addition. It does not matter what the distance is between the different elements that need to be added together, the time taken

to implement the addition is the same (in clock cycles). Communication can be effected concurrently which leads to further speed ups and allows the divide and conquer algorithm to be applied effectively and efficiently.

- 5 There is a particular benefit when both the alternation network and the communications network are combined together, the resultant speed up of operation provides a significant advance over the prior art.

10 According to a further aspect of the present invention there is provided a combination of an alternation network as described above and a content addressable memory, the content addressable memory being arranged to receive the selected one of the odd and even match vectors for use as a matching result vector in a further iteration of the divide and conquer algorithm.

15

This combination advantageously leads to a compact and relatively fast improved hardware implementation of the desired divide and conquer algorithm.

- 20 Furthermore the present invention also extends to a combination of a communication network as described above, a content addressable memory, and a plurality of addition registers.

25 In order that the invention may be more readily understood, some embodiments in accordance therewith will now be described, by way of example, with reference to the accompanying drawings, in which:-

Figure 1 is a schematic example showing the 'divide and conquer' principle;

- 30 Figure 2 is a circuit diagram showing a hardware implementation of the basic functional logic block of the alternation network;

Figure 3 is a circuit diagram showing a novel implementation of the alternation network of Figure 2;

Figure 4 is a circuit diagram showing the inject network of Figure 3 in greater detail;

- 5 Figure 5 is a circuit diagram showing an alternative to the network shown in Figure 3;

Figure 6 shows a modular network comprising multiple segments of the network as shown in Figure 3;

10

Figure 7 shows a modular network comprising multiple blocks of the network as shown in Figure 6;

- 15 Figure 8 shows a modular network comprising multiple chips implementing the network as shown in Figure 7;

Figure 9 is a schematic of a network according to the invention describing in detail the role of a forward-fed correction term;

- 20 Figure 10 shows an efficient circuit design catering for redundancy in the event of a faulty microprocessor component;

Figure 11 is a flow chart covering the steps in carrying out a 'divide and conquer' operation in accordance with the present invention;

25

Figure 12 shows a data propagation network for movement of data from a source to a destination;

Figure 13 shows the data propagation network of Figure 12 in more detail;

30

Figure 14 shows a bypass circuit compatible with the data movement network as shown in Figures 12 and 13;

Figure 15 shows how bypass features may be enabled at a higher level to that shown in Figures 13 and 14; and

Figure 16 shows a network for facilitating synchronisation between multiple microprocessor chips.

By way of background, a typical SIMD processor with which the invention can be reduced to practice is the Associative String Processor (ASP) data parallel processor manufactured by Aspex Technology Limited. The ASP is a SIMD data processor that in typical configurations operates on 1024 to 65536 data items in parallel. One of the features of the ASP processor is the string (one-dimensional) topology of the processor array, wherein the processing elements are connected end-to-end. This facilitates some of the features of the network described in the following, although the principles of the invention may obviously be applied to other processors, SIMD or not. It is noted that the ASP is associated to a portion of memory that is referred to as 'local memory', and the ASP may be programmed by a Central Processing Unit (CPU) which is able to off-load tasks to the ASP, for example processing data stored in the local memory associated to that ASP. This locally accessible memory is referred to as Content Addressable Memory (CAM).

An improved version of the alternation network architecture according to an embodiment of the present invention is illustrated in Figure 3 which shows a modular alternation network.

25

The module 6 labeled 'alternation network' corresponds to the basic alternation network shown in Figure 2. Typically this portion of the network comprises only a few stages (processing elements), since propagation through this network must be carried out twice, the first time to allow propagation to occur throughout the alternation network, reaching a static output and determining a correction term. The second time is when applying the correction term to produce the final even-odd sets. Consequently this sub-network must be kept as short as possible, consistent with the requirements of the floorplan.

It is helpful here to note that the alternation network 6 as shown may comprise a number of structures as shown in Figure 2. Typically 16 ASP processors arranged in this way in a network are known as a 'segment'. Following the segment, four segments arranged together are known as a 'group', four groups arranged together known as 'block', four blocks arranged together known as a 'quad' and four quads arranged together to form a 'chip'. This hierarchical arrangement can be thought of as a microprocessor chip having various subnets of smaller groups of standardised sub-components.

In Figure 3 there are sixteen ASP processors connected in the alternation network 6, and this building block shall be referred to as a 'segment' hereafter.

The module 8 labeled 'inject network' determines what the initial input to the alternation network based upon a number of factors, or it may be an assumed value. One of the values that may comprise the input to the alternation network is obtained from the output of an EXOR gate 10, which may feed through to the alternation network a correction term from a neighboring alternation network or a correction term from a network outside of the present network, for instance a block bypass, in order to speed up the overall processing of the microprocessor chip. The speeding-up function will also be described in greater detail later.

As several copies of the arrangement of Figure 3 may be connected together, AND gate 12 serves as a partition switch whose function is controlled by a signal P_s where s denotes the segment identifier. For example $s=0$ is considered to be the leftmost segment, the starting segment with no segments to its left. If this signal is TRUE then the carry forward of the correction term (LAI1) to that segment via the partition switch to that alternation network segment is disabled. The correction term LAI1 is otherwise derived from a neighbouring segment.

If the partition signal is TRUE then communication into this segment is also disabled, and the inject network 8, which controls the insertion of a

propagating signal (LAI1) into the leftmost end of the modular alternation network, as well as responding to the partition flag (P_s) to regulate input to the network.

- 5 The inject network is shown in detail in Figure 4 and comprises three multiplexers 16, 18 and 22 and an isolation flip-flop 20. Multiplexer 22 responds to the state of a partition flag P_s for the segment. If the term P_s is true it means that particular segment is at the start, or is the left-most segment, and so no correction term will ever be applied to it. Accordingly,
- 10 when P_s is true the inject condition 'I' will be applied to start processing in that segment. Inject condition I can be programmable true or false. If P_s is false then a correction term will be applied at a later step when the correction term is ready.
- 15 Apart from being used to determine whether this segment is partitioned from the segment to the left, the inject condition I can also be used to determine whether the leftmost participating processor in this segment belongs to the even or odd set. This can be determined from the number of results in the sets of odd/even results in the preceding segment. If the partition flag P_s is
- 20 TRUE, this segment is partitioned from segments to its left. The global inject condition then overrides all forward propagating correction terms into the network input and unconditionally becomes the input to the ALTERNATION NETWORK.
- 25 The multiplexer 16 at the far left of the inject network is also dependent upon the partition flag P_s . If P_s is false then the initial state (which is frequently tied to ground or '0') is passed through to multiplexer 18 which responds to a correction signal C. The correction signal C determines whether the inject network continues to loop (until C becomes false) or whether the correction
- 30 signal determined by the alternation network is ready to be injected back into the alternation network (once the correction term has been determined and is stable). The correction signal C is only inserted back into the alternation network when an initial static result of odd/even sets has been determined based upon the initial assumed matching value. When C is false the

correction signal is clocked through the isolation flip-flop which also flips to allow the correction term to be injected into the lowest order network and apply any correction as appropriate to the set of odd/even results. The corrected set of odd/even results is stored as vectors, and the way these
5 vectors are stored is described in greater detail below.

It is remembered that upon initiation of an odd/even alternation process the initial starting value against which M_n values are compared (which represents the total number of odd and even matches in the preceding segment(s)) is
10 assumed – normally '0' – and once the initial set of odd/even results are generated based upon this initial assumed starting value, the correction term is then applied to ensure that the odd/even alternation in this segment follows from the odd/even split in the previous segment(s) – the correction term can in this sense be considered as a parity check. Obviously a correction step is not
15 needed if this segment is the first in line, which is when the inject condition is used instead.

In operation, there are two distinct phases – the first phase is the generation phase when the set of odd/even results are being propagated through the
20 alternation network based upon the initial assumed matching value. Whilst this is going on the flip-flop prevents the correction signal (determined from a preceding segment) to be inserted into the current segment and the correction signal C is true. Once the initial set of odd/even results have been generated, the correction signal C stops, at which point the correction term is applied
25 through the isolation flip-flop and allows the correction term to be applied to the lowest level network. Any corrections are then propagated through the alternation network to arrive at the final, corrected set of odd/even results.

In other words, the communication network will be generating a carry forward
30 correction term into the alternation network input which will be used to modify the initial state of the alternation network (provided that P_s is FALSE). The initial state of the alternation network is used as the basis of generating the correction terms and propagating them forward through the network, hence

the correction term is isolated from the alternation network by the isolation flip-flop.

Once the correction term has been generated and is stable, it will be clocked
5 through to the alternation network during the second phase of alternation network operation, when the correction term will be enabled through the isolation flip-flop and applied to the initially-determined set of odd/even results.

When the signal C is false, then the isolation flip-flop will be pre(set) or
10 pre(reset) depending upon the state of the inject condition I and the partition flag P_s .

In an alternative embodiment for deriving a set of odd/even results, a second alternation network may be provided to run the correction term at the same
15 time as the first alternation network is determining the propagation signal and correction term for the following segment. An example of this arrangement is shown in Figure 5.

In Figure 5 there are two sets of alternation networks Alt1 and Alt2 both
20 receiving the same set of M_n data for comparison. Alt1 is connected at the starting end with a multiplexer 26 similar to those shown in the inject network of Figure 4, except that the partition flag P_s determines whether the initial inject condition I is used (if the segment was the first after a partition), or whether '0' is to be used as the initial assumed value. While propagation takes
25 place through alternation network Alt1 to determine the set of odd/even results based upon the assumed value I, at the same time the lower alternation network Alt2 is determining the same set of odd/even results using the same set of M_n data inputs as Alt1 is using. However, the alternation network Alt2 is also receiving a forward propagated correction term via EXOR gate 10 which,
30 itself will have been generated by a preceding alternation network. As a result of this network's asynchronous operation, segments further away from the start may have several correction terms applied to its set of odd/even results until a static condition is achieved.

The set of odd/even results generated by the alternation network Alt2 is stored in a set of registers, referred to as TAG registers, for further use.

- 5 The network shown in Figures 2 to 5 achieves the functionality of generating alternating match sets of odd and even items at a vastly improved speed over the current art, and as mentioned before this is required for the common task of 'divide and conquer' calculations. This improvement in speed is achieved by generating an initial alternate odd and even appropriation locally within a
10 segment, then generating a correction term based on the forward propagation of the output of the local EXOR chains from preceding segments. The correction term, when it has reached a static value, provides information about whether the total number of matches within a segment are even or odd. This correction term is forward propagated through the network, utilising a
15 hierarchy of bypass levels (which are described in greater detail below), which further accelerate the propagation of correction terms. In a subsequent cycle, the correction term is utilised to generate a new input to the local segment, which then serves to cause the segment to maintain or reverse initial alternate order to achieve the correct final state.

20

The module shown in Figure 3 can now be readily cascaded to form a higher-order network. For example, Figure 6 shows a number of identical alternation network segments of the type shown in Figure 3 cascaded to form a block network.

25

- The block partition switch 30 is controlled by a block partition flag P_b which is the logical combination of individual segment partition flag P_s control signals. In this way, only if all segment partition switches 12 (not shown in Fig 6) within the span of a block are closed is the forward propagation of a correction term
30 enabled. This is to prevent alternation networks within the block being bypassed which have been tasked to process relevant data.

The hierarchical arrangement described above can be extended indefinitely, for example by the cascading block-level networks as shown in Figure 6,

accompanied by the addition of the corresponding external network as shown in Figure 7.

In the example of Figure 7 it is assumed that a chip comprises an array of blocks, as per the definitions above. The block bypass partition switch 30 and the microprocessor chip partition switch 34 perform similar functions to the segment bypass partition switch 12 and the block partition switch 30 of Figure 6. For clarity, microprocessor chip will hereafter be referred to simply as 'chip'.

It is anticipated that this hierarchical network topology can be readily extended limited only by the constraints of the target speed of the correction phase and the floorplan topology of the integrated circuit. It is anticipated therefore that the hierarchy within a given chip may be extended to comprise more than the two levels described above to comprise any number of levels or subnets.

Ultimately however, the communications signals must propagate off-chip, supporting the goal of transparent chip-chip communication, as well as supporting the continuation of the network logical hierarchy with a minimum of external logic.

Figure 8 illustrates how the network hierarchy may be readily extended across multiple chips on a board. In this example a signal EBE has been introduced which is assumed to be TRUE if the chip may be bypassed, that is no open switch exists anywhere within that chip on any hierarchical level. Board bypass switches 38 are set true or false according to the signal EBE.

The examples shown in Figures 3 to 8 demonstrate the readily-expandable and modular nature of this type of chip design, which is extremely flexible and can be tailored to suit the demands of the application.

The operation of a hierarchical network as described above is illustrated in Figure 9, which shows a representative block comprising three segments. In this example, as no partition flags are set the network is continuous and unbroken along its length. The initial state of the individual segments is

evaluated in parallel (assuming a default inject condition set FALSE), resulting in sets of even and odd match counts. These are then combined by the carry forward network as described in Figures 3 to 5 to produce the correction terms. The role of the EXOR networks becomes clear in this illustration since

5 the carry forward of pairs of odd correction terms (FALSE, FALSE) or pairs of even correction terms (TRUE, TRUE) always results in the generation and propagation of an even (TRUE) carry forward term. Thus, the alternating odd/even continuity is maintained across the three segments. Similarly, the application of, say, an even and an odd correction term (TRUE, FALSE)

10 results in the propagation of an odd correction term (FALSE) to maintain odd/even continuity in a following segment.

The duration of generation and propagation of the correction terms is determined by the internal modular hierarchy within the network. The final

15 stable correction terms are applied on a subsequent clock edge, when the isolation flip-flop in the inject network is clocked – this is described in the narrative accompanying Figure 4. Following the application of the correction term, a sufficient interval must be allowed for the correction term to propagate through the individual segment networks to achieve the final even-odd split.

20 This interval will be designed to be less than one clock cycle.

Whenever an operation is extended (such as the generation and propagation of the correction terms) this is implemented in the context of a synchronous processor device by issuing repeated instructions according to a deterministic

25 (compiled) algorithm.

It is found that a typical usage of the network requires a sequence of one or more of the correction data processor instructions, followed by an evaluation processor instruction and a single propagation instruction. Control logic within

30 the chip instruction unit allows the length of this sequence to be deterministic.

Example I

Consider the following example. The instruction sequence that needs to be fed to the data processor for a particular application is:

```
{{A}RepeatSlot(4), B }
```

5

where {...} RepeatSlot(N) means repeat the instruction inside the braces N times.

Example ASP alternation Code

10

Consider the following example. Count participating APEs and write the count into the data register in the rightmost APE in the string.

The pseudo-code for this is:

15

```
identify participating APEs;
A = 0;
FOR X = LSB TO MSB DO
    alternate participating APEs into even and odd sets;
20 IF rightmost APEs is in odd set THEN
    A[rightmost][X]=1; (* Xth bit of rightmost APE is
    set *)
    mask odds and reassign evens as new participating APEs
END
```

25

The actual ASP code to perform this task is given below:

Initialise

```
{
30 {InitWriteSerial(every,sloOpen,dX,dX,0,0,aw{ab2.d0,ab1.d0}
,NoRead)},
    {ExecWrite(bmD,NoCsoOpt)}};
```

Start of loop, initialise for counter and set intra and

```
35 inter chip inject inputs to TRUE
    {Set(agbInject),
      For(sfsz,1),
        {SegmentInject(1)},
        {NoOp()}};
```

40

Finding participating processors and generate alternate even-odd subsets

```
{
{InitSearchOrAddSerial(dX,dX,0,0,aw{ab2.d0,ab0.d1}}),
```

```

        {ExecSearchOrAddSerial (bmD, trAlt, NotForceZero),
          RepeatSlot (AlternatesSlots) }};

Set intra-chip inject inputs FALSE
5   { {SegmentInject(0)},
      {NoOp()}};

Set inter-chip inject input FALSE
10  {Clear(agbInject),

     {InitWriteSerial (at, s1Open, dX, dX, 0, 0, aw{ab1.d1}, NoRead) },
     {ExecWrite (bmD, NoCsoOpt) }};

Establish that rightmost participating APE is odd
15  {
     {InitWriteTernary (an1, s1Open, tf0, XT, ScalarBus, aw{ }) },
     {NoOp() }};

     {
20  {InitWriteSerial (a11, s1Open, dX, dX, 0, 0, aw{ab2.d1}, NoRead),
      RepeatSlot (ActivateSlots) },
      {ExecWrite (bmD, NoCsoOpt) }};

     {
25  {InitSearchOrAddSerial (dX, dX, 0, 0, aw{ab0.d1, ab1.d0, ab2.d0}
        ) },
        {ExecSearchOrAddSerial (bmD, tr1, NotForceZero) }};

     {
     {InitWriteTernary (anr, s1Open, tf0, XT, ScalarBus, aw{ }) },
     {NoOp() }};

30  Write result LSB
     {Init (sfal, sfbal),

        {InitWriteSerial (aerf, s1Open, d1, dX, sfal, 0, aw{ab2.d1}, NoRe
35  ad),
        RepeatSlot (ActivateSlots) },
        { },
        Inc (sfal) };
        { { },
          {ExecWrite (bmD, NoCsoOpt) }};

40  Mask current odd set from participating in next cycle of
alternates
     {
     {InitWriteSerial (every, s1Open, dX, dX, 0, 0, aw{ab1.d0}, NoRead
45  ) },
        {ExecWrite (bmD, NoCsoOpt) },
        EndFor () };

```

The evaluation of this task using the invention would result in completion of this application in 240 clock cycles instead of 304 clock cycles using the network shown in Figure 2, which represents the prior art. The effect is far more significant in the next planned ASP device, which will incorporate 4096 APEs. Here completion of the algorithm will require some 400 clock cycles, compared to some 2500 clock cycles without the invention.

The hierarchical nature of the network topology allows a subnet (at segment, block, chip etc. levels) to be readily bypassed in the event of a defect which results in faulty operation. For example, Figure 10 shows the addition of an isolation gate 44 in the output of an alternation network 46 forces the output to a null state in the event of a fault in the alternation network 46. By forcing the partition switch P unconditionally closed existing hardware may be exploited to allow any forward propagating correction terms to unconditionally bypass the segment, block or chip etc. Whilst built-in redundancy is a known technique, this chip design enables a particularly space efficient method of achieving redundancy at all levels of the hierarchy.

In practice, an implementation of the invention can be found in the VASP-4096/TX microprocessor. This microprocessor design completes a full alternation across 4096 data parallel processor in a simulated interval of 24ns - or 5 clock cycles at 160MHz. It is implemented in a single 0.18µm silicon process.

In the above embodiments, once an alternating odd and even set has been identified it is usual for the values identified by the odd and even sets to be summed. It is worth remembering that only memory addresses that have been successfully matched with the match signal are alternated odd and even, and once the alternation has taken place, the values contained at those memory addresses need to be processed, usually by summation. In previous designs summation takes place wherein the value of each memory address is checked against the match signal, and if true, the value contained in that memory address is added to a previous sum and the next memory address checked in

turn. This synchronous summation is time consuming, especially in a case where adjacent memory addresses are far apart in memory address space.

5 The hardware embodiments shown in Figures 2 to 10 demonstrate that the identification and alternating actions of the overall summation procedure can be achieved through chip design alone. Furthermore, through use of bypass networks at all levels of the modular hierarchy, the speed of the chip may be increased significantly. For a large-scale 'divide and conquer' procedure involving many thousands of data items, many iterations of the 'identify and
10 alternate' process need to be carried out before all these data items have been summed together. An overview of this procedure is provided in the flow chart of Figure 11.

Starting at 50, the CPU sends instructions to LAC controllers to identify
15 required data items. Next, at 52 the LACs determine the data items from their associated CAM memory that match the required parameters, and by processing these data items through alternation networks at 54, these data items are split into odd and even pairs and the resulting sets of odd and even results are written as vectors into TAG registers at 56. Remembering that
20 each data item corresponds to a data register containing a binary word of data, a movement of data then takes place at 58 of the data words from the data register of data 'sources' (which can be user-defined as corresponding to 'odd' or 'even' data items) to the data register of data 'destinations', which will be the 'odd' or 'even' complement of the data sources. Register additions may
25 then take place and the results written to the destination data register at 60. A check is then carried out at 62 to see if there is more than one data item in the 'even' TAG register. If there is, the data in the 'even' TAG register is written to the CAM at 64 and the procedure loops back to 54 to continue splitting data items into 'eve' and 'odd' pairs until the final check at 62 determines that there
30 is only one data item left. When there is only one data item left, the check at 62 returns a negative finding in which case the data in the remaining data register is the answer to the 'divide and conquer' summation.

As can be seen from the above, the movement of data items from an 'even' data register to an 'odd' data register is crucial to the overall speed of the network in carrying out the 'divide and conquer' summation. It is easily imaginable that after several splitting and summation steps that the remaining data items left are spaced far apart. The usual procedure for data movement is synchronous data movement, whereby the contents of the data register corresponding to the source processor is shifted one processor along until it reaches the destination processor's data register. Each shift will take a clock cycle so synchronous movement will be very slow compared to the asynchronous data movement achieved in the alternation network stages.

An improved summation network referred to as a 'remote activation network' is disclosed in Figure 12 wherein asynchronous summation is enabled in hardware. From step 56 of Figure 11 wherein the 'even' and 'odd' data items are written to respective TAG registers, the data items in the 'even' TAG register are bound as 'sources' and the data items in the 'odd' TAG register are bound as destinations. This binding convention may be user defined so that the 'odd' TAG register data items are considered as 'sources'. The idea here is that every source data item will aim to have its data transported to its adjacent destination data item as quickly and efficiently as possible.

Turning to Figure 12, a single bit of data $TxDATA_0[bit]$, for example which is HIGH, from the first source data item's data register is added with a source data item flag S_0 (which is set as TRUE) and added into AND gate 70. The output of AND gate 70 is then entered into an OR gate 72 with the inverse of the source flag S_0 ORRED with destination flag D_0 . The destination flag D_0 is only TRUE at location of the destination data item so D_0 is LO, and the result of OR gate 72 is TRUE as S_0 is HIGH. The output of OR gate 72 is fed into AND gate 74 which has as its other input the output from multiplexer 76, which nominally has a TRUE output. Thus the output of AND gate 74 is true, and to the right of the combination of logic gates 70, 72 and 74 are other logic gates arranged identical to logic gates 70, 72 and 74. As all these other logic gate arrangements do not correspond to a source or destination data item, the TRUE output from AND gate 74 is propagated along the series of AND gates

(not shown in Figure 12) until the destination data item is reached. At the destination data item, AND gate 80 outputs a FALSE value as $TxDATA_n[bit]$ and S_n are both zero. However, as S_n ORRED with D_n is TRUE, as D_n is HIGH, the inverse of this is LO, so the output of OR gate 82 is FALSE, which means AND gate 84 does not give a TRUE output and stops the propagation of the TRUE output from AND gate 74 along the chain of AND gates. Finally, AND gate 88 provides a TRUE output as both the output from AND gate 74 and D_n are HIGH. The resultant output from AND gate 88 is written into the $RxDATA_n[bit]$ destination data register as a HIGH for that data bit. The same works in the case of the source data bit being LO to begin with – in this case rather than AND gate 74 propagating a TRUE output from left to right in Figure 12, a FALSE output would be propagated instead. This FALSE output would still be stopped by the destination AND gate 84 due to the condition of S_n and D_n , and the resultant FALSE value from AND gate 88 would be written to $RxDATA_n[bit]$ destination data register as a LO value for that data bit.

Once all the data bits from the source data register have been bit-serially transferred to the destination data register, a simple register addition action can be accomplished to sum the transferred data with the existing data in the destination data register to produce the new value which is written to the destination data register.

Figure 13 shows the example of Figure 12 in more detail, in this case the destination data item with corresponding logic gates 80, 82, 84 and 88 is two processors away from the source data item, with corresponding logic gates 70, 72, 74 and 78. Intermediate processors, which do not have an 'odd' or an 'even' data item, and which have corresponding logic gates 90, 92, 94 and 98, do not affect the propagation of the $TxDATA_0[bit]$ from the output of AND gate 74 to the AND gate 88 of the destination data item.

In other words, a communication signal is generated from a source node (denoted by $S=TRUE$) that is then propagated in the direction towards the right side of Figures 12 and 13. This data is generally the logical product of a source 'flag' (i.e. S) and a data bit. This data is then propagated transparently

through intermediate processors that are neither a source nor a destination (i.e. S=FALSE and D=FALSE), terminating at a processor that is marked as a destination (i.e. D=TRUE).

- 5 It is obviously desirable that in order to achieve the maximum processing speed, the data movement network as shown in Figures 12 and 13 is able to move data as fast as possible, especially considering the speed of operation of the alternation network that preceded the data movement network. It is noted that the transfer of data between source and destination nodes, as
- 10 described in the narrative to Figures 12 and 13, operate essentially as circuit switched data between a source and a destination, and as no two circuits will ever overlap (due to the alternating 'odd' and 'even' split), data movements between each pair of source and destination can take place simultaneously.
- 15 Furthermore, as in the alternation networks, bypass switches may be employed so that segments which do not have a part to play in a circuit between a source and a destination may be bypassed entirely, leading to further speed improvements. Figure 14 shows this with bypass switch 100. Obviously it is only appropriate to bypass whole segments in cases where
- 20 bypassing does not progress past isolation switches that occur periodically in the network. One such topology might have such switches every segment. The External Bypass Enable flag which sets the segment bypass switch 100 only allows segment bypass if the following two conditions are met: (1) no isolation switch is open within the span (i.e. all closed); and (2) no destinations
- 25 are present in the span to be bypassed.

A further variant of the bypass function is shown in Figure 15 wherein three segments of an asynchronous data movement network exemplified in Figures 12 to 14 are able to be bypassed at once if the EBE bypass AND gate 102 is

30 TRUE.

As mentioned before, it is preferable that this bypass topology be repeated hierarchically. The topologies that have been reduced to practise have

generally adopted a bypass granularity in multiples of four, i.e. bypassed processors grouped in x16, x64, x256, x1024, x4096 etc.

5 It is envisaged that allowing propagation of communicated data to occur – and complete – within a single clock cycle is preferable. It is recalled that the operation of this network – although asynchronous – occurs within a synchronous operating framework, i.e. the beginning and end of communication cycles are bounded by clock cycles. Failure to complete the propagation in a single clock cycle will generally require multiple clock cycles
10 to be allocated – defining the communications path a multi-cycle path.

Moreover in the context of a system comprising communicating processor chips (i.e. extending this network across multiple devices), it is unlikely that the operations between chips can be easily maintained as clock synchronous.
15 Therefore a method of achieving clock level synchronism – via an appropriate mechanism – is achieved in order to make this invention viable as a scaleable system. Such a synchronisation scheme involves a software programmed *wait()* statement which is linked to a hardware rendezvous. Such a rendezvous may be signalled by multiple generators – which must all reach
20 the synchronisation point before the communications operation may proceed. The synchronisation is achieved using a pair of rendezvous signals, *ape Rendezvous Request* (*apeRDVR*) and *ape Rendezvous Complete* (*apeRDVC*). This synchronisation may occur between autonomous chips. Figure 16 illustrates this.

25

In Figure 16, all the *apeRDVR* signals from the devices are combined using an AND gate, ensuring that all chips are ready to communicate. The resulting signal is optionally resynchronised to the master chip clock - the output of which is connected to all the chip *apeRDVC* inputs to signal that the
30 communication may proceed.

Naturally, although the foregoing descriptions suggest a communication topology that routes from left to right, in fact such a network is beneficially extended to include operation which is simplex (i.e. configurable to

communicate leftwards or rightwards) and duplex (i.e. duplicate networks are implemented, allowing simultaneous communications leftwards and rightwards).

- 5 Having described particular preferred embodiments of the present invention, it is to be appreciated that the embodiments in question are exemplary only and that variations and modifications such as will occur to those possessed of the appropriate knowledge and skills may be made without departure from the spirit and scope of the invention as set forth in the appended claims.

Claims

1. An alternation network for use with a content addressable memory for
5 implementing a divide and conquer algorithm, the network comprising:
a plurality of alternation modules connected in series together, each
module comprising:
a plurality of cascaded logic gates arranged to propagate a
match parity signal via the gates along at least part of a matching result
10 vector, the matching result vector being generated by execution of a
matching instruction on the content addressable memory, and the logic
gates being configured to change the parity of the match parity signal in
accordance with the matching result vector; and
a vector output arranged to output a parity level vector of the
15 propagated match parity signal present at the each gate of the plurality
of logic gates;
a logic network for dividing the matching result vector into an odd
match vector and an even match vector representing respectively odd and
20 even numbered elements of the matching result vector, by use of the parity
level vector; and
means for writing a selected one of the odd and even match vectors to
the content addressable memory.
- 25 2. An alternation network of Claim 1, wherein the plurality of cascaded
logic gates are arranged to implement a cascading exclusive-OR function on
at least part of the matching result vector.
3. An alternation network of Claim 1 or 2, wherein each module further
30 comprises a parity output for outputting a resultant parity of the plurality of
cascaded logic gates, after the match parity signal has propagated through all
of the logic gates of the module.

4. An alternation network of any preceding claim, further comprising generating means for generating a feed-forward term, the generating means being arranged to utilise the parity output of a neighbouring alternation module.
- 5
5. An alternation network of Claim 4, wherein the generating means is arranged to utilise a correction term, the correction term being generated from a non-neighbouring alternation module.
- 10
6. An alternation network of Claim 5, wherein the correction term represents the parity of a part of the matching result vector as determined by a section of the alternation network comprising a plurality of alternation modules.
- 15
7. An alternation network of any of Claims 4 to 6, wherein the vector output is arranged to be generated after a received feed-forward term has been applied to the module.
- 20
8. An alternation network of any of Claims 4 to 7, further comprising a plurality of interconnection pathways, means for determining a correction term from the parity of combined feed-forward terms, and pathways for feeding back the correction term to one or more alternation modules.
- 25
9. An alternation network of Claim 8, wherein the correction term is generated from a hierarchy of interconnection pathways.
- 30
10. An alternation network of any of Claims 4 to 9, further comprising a partition switch for each alternation module, the partition switch being arranged to disable propagation of the parity of the feed-forward term to any other alternation module.
11. An alternation network of any preceding claim, wherein the plurality of cascaded logic gates comprises a first cascade of logic gates arranged to calculate the parity of the matching result vector; and

a second cascade of logic gates for implementing an exclusive-OR function on the matching result vector using the received feed-forward term.

12. An alternation network of Claim 11, wherein the first cascade of logic gates is arranged to implement an exclusive-OR function on the matching result vector using a predetermined logic level input to a head of the plurality of cascaded logic gates.

13. An alternation network of any of Claims 1 to 10, wherein each module further comprises an injection module for receiving a starting logic level for the module and injecting the same into the plurality of cascaded logic gates.

14. An alternation network of Claim 13, wherein the injection module is arranged:

15 to assume an input logic level;
 to inject this into the plurality of cascaded logic gates;
 to compare the assumed input logic level with the received starting logic level and
 to inject the starting logic level into the plurality of cascaded logic gates
20 if there is a difference.

15. An alternation network of Claim 13 or 14, wherein the parity output is arranged to create the parity level vector from the injection of the assumed input logic level when the assumed input is correct and to create the parity level vector from the injection of the starting logic level when the assumed input is incorrect.

16. An alternation network of any preceding claim, wherein the logic network comprises an AND gate for generating an element of the even match vector from the parity match signal and the matching result vector.

17. An alternation network of any preceding claim, wherein the logic network comprises an inverter for inverting the propagated parity match signal

and an AND gate for generating an element of the odd match vector from the inverted propagated parity match signal and the matching result vector.

18. An alternation network of any preceding claim, wherein the logic
5 network comprises means for storing the odd and even match vectors.

19. A combination of an alternation network as claimed in any preceding
claim and a content addressable memory, the content addressable memory
being arranged to receive the selected one of the odd and even match vectors
10 for use as a matching result vector in a further iteration of the divide and
conquer algorithm.

20. A combination according to Claim 19, further comprising a plurality of
addition registers, each addition register being associated with an element of
15 the content addressable memory.

21. A combination according to Claim 19 or 20, wherein the content
addressable memory comprises a plurality of strings of array processors.

20 22. A communications network for use with a content addressable memory
for moving data between remote locations, the network comprising:

dividing means for dividing a matching result vector generated by
execution of a matching instruction on the content addressable memory, into
an odd match vector and an even match vector representing respectively odd
25 and even numbered elements of the matching result vector;

assigning means for assigning the odd and even match vectors as
source and destination locations;

a plurality of cascaded logic gates arranged to propagate data via the
gates along at least part of the matching result vector;

30 input means for inputting a bit of a stored data value associated with a
source location into the plurality of cascaded logic gates;

applying means for applying the source and destination results to the
plurality of cascaded logic gates to propagate the bit of stored data at the
source location to the destination location; and

extracting means for extracting the input bit from the plurality of cascaded logic gates at the destination location.

23. A communications network of Claim 22, wherein the network is
5 arranged to propagate a bit of the data stored at each source location to a corresponding destination location concurrently.

24. A communications network of Claim 22 or 23, wherein the plurality of cascaded logic gates is arranged to implement a cascading AND function.

10

25. A communications network of any of Claims 22 to 24, wherein the plurality of cascaded logic gates comprises a plurality of groups of sequentially connected logic gates.

15 26. A communications network of Claim 25, wherein each of the plurality of groups comprises an activation switch for activating or deactivating the respective group of sequentially connected logic gates.

27. A communications network of any of Claims 22 to 26, further
20 comprising bypass means for connecting one group of sequentially connected logic gates to another non-neighbouring group of sequentially connected logic gates, thereby bypassing a neighbouring group of sequentially connected logic gates.

25 28. A communications network of Claim 27, wherein the bypass means is arranged to be activated when the neighbouring group of sequentially connected logic gates does not comprise any destination location.

29. A communications network of Claims 27 or 28, wherein the bypass
30 means is arranged to be activated when the neighbouring group of sequentially connected logic gates is deactivated.

30. A communications network of any of Claims 27 to 29, wherein the bypass means is arranged to be activated when a plurality of neighbouring

groups of sequentially connected logic gates do not comprise any destination location.

31. A communications network of any of Claims 27 to 30, wherein the
5 bypass means is arranged to be activated when a plurality of the neighbouring groups of sequentially connected logic gates are deactivated.

32. A combination of a communication network as described in any of
10 Claims 22 to 31, a content addressable memory, and a plurality of addition registers.

33. A combination of Claim 32, wherein the combination is arranged to
implement a remote addition algorithm, whereby data stored in the plurality of
addition registers of selected source and destination locations are combined
15 together using the communication network.

34. A combination of Claim 33, wherein the content addressable memory is
arranged to receive the selected one of the odd and even match vectors for
use as a matching result vector in a further iteration of the remote addition
20 algorithm.

35. A combination of any of Claims 32 to 34, wherein the dividing means
comprises an alternation network as claimed in any of Claims 1 to 18.

25 36. A combination according to any of Claims 32 to 35, further comprising a plurality of addition registers, each addition register being associated with an element of the content addressable memory.

37. A combination according to any of Claims 32 to 36, wherein the content
30 address memory comprises a plurality of strings of array processors.

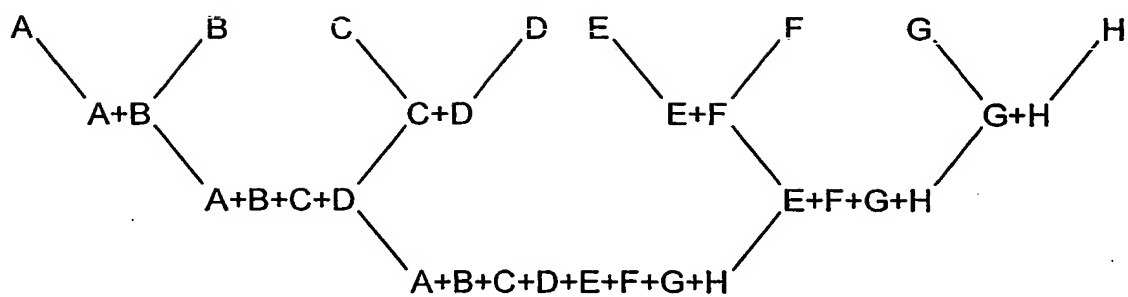


FIG. 1

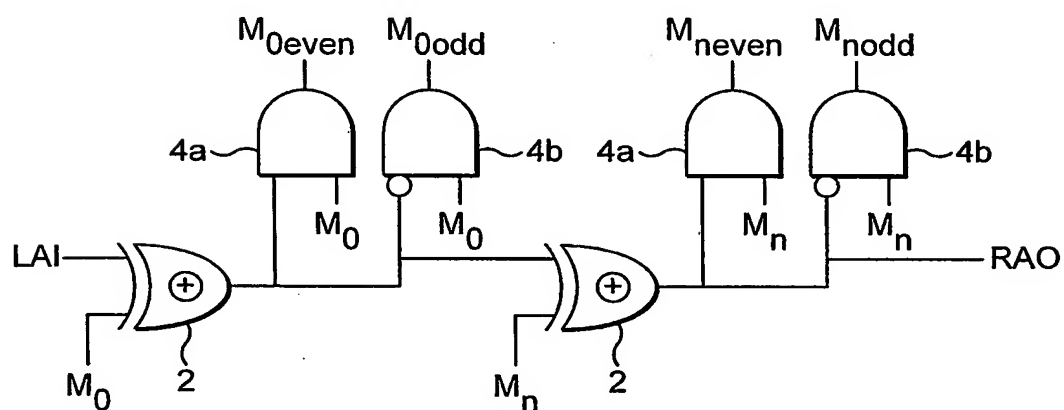


FIG. 2

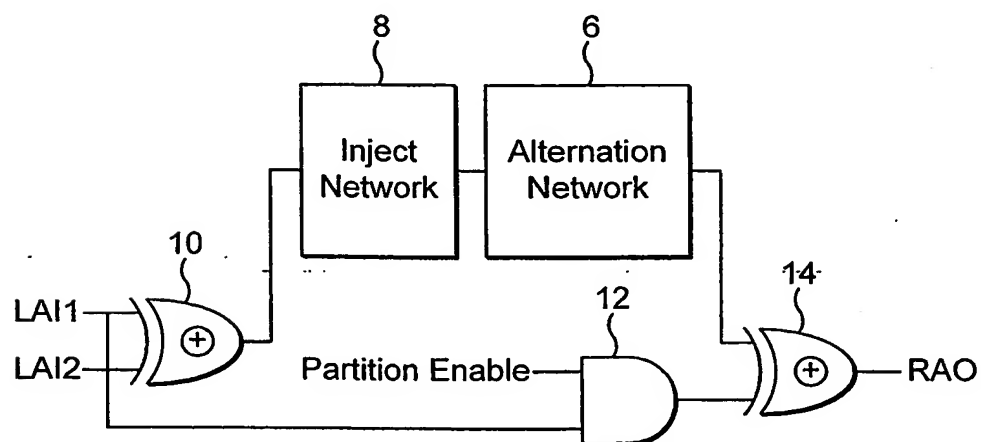


FIG. 3

2 / 8

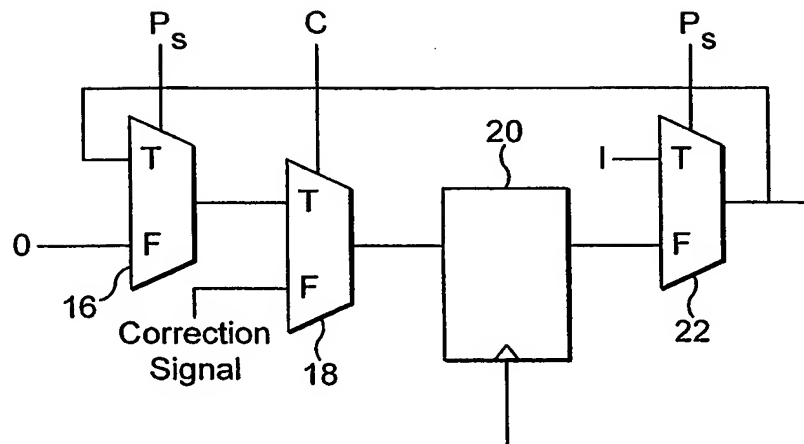


FIG. 4

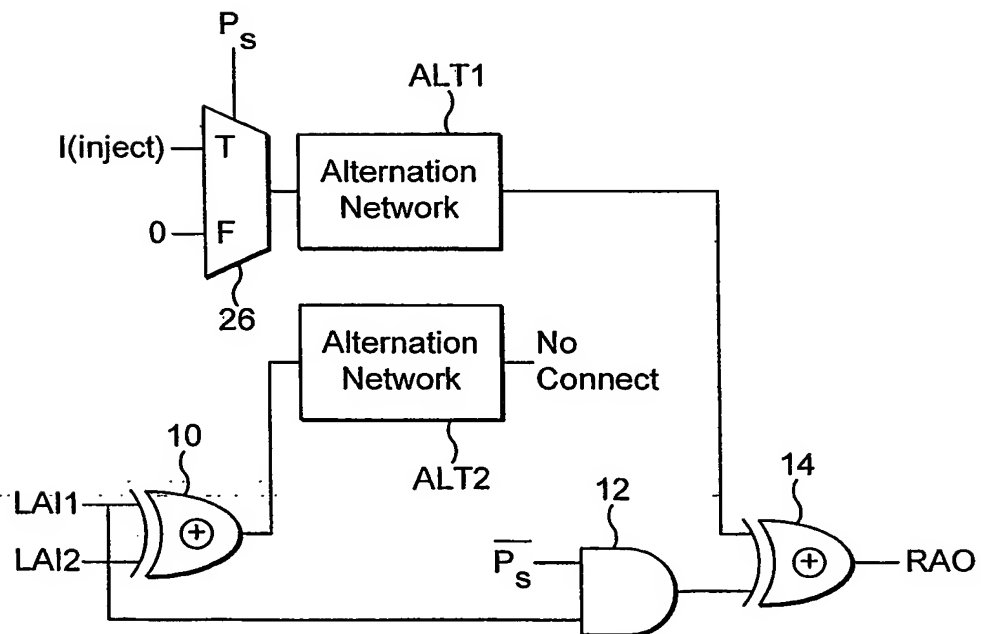


FIG. 5

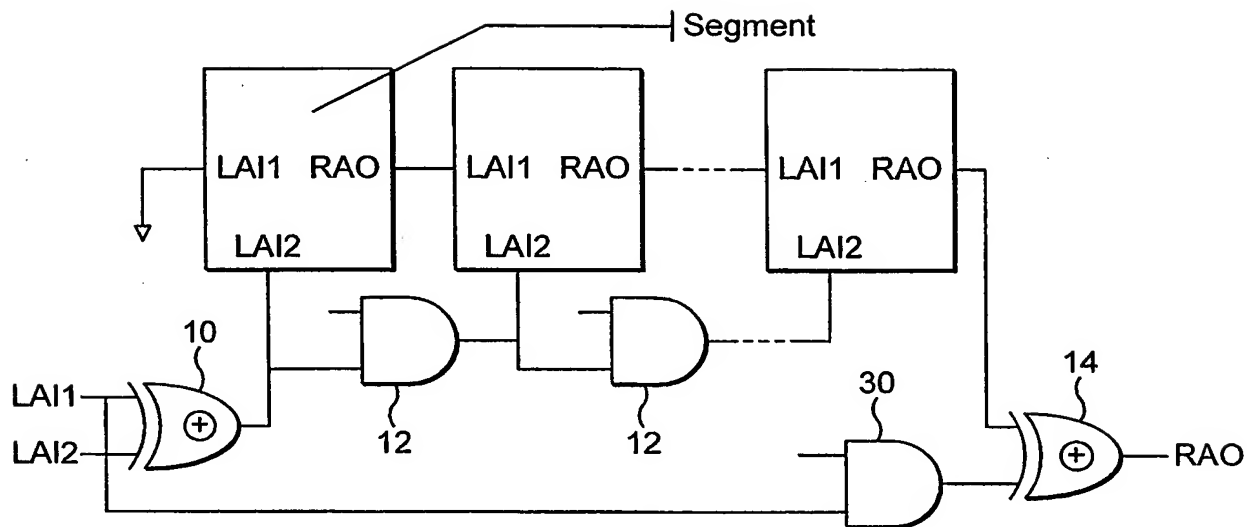


FIG. 6

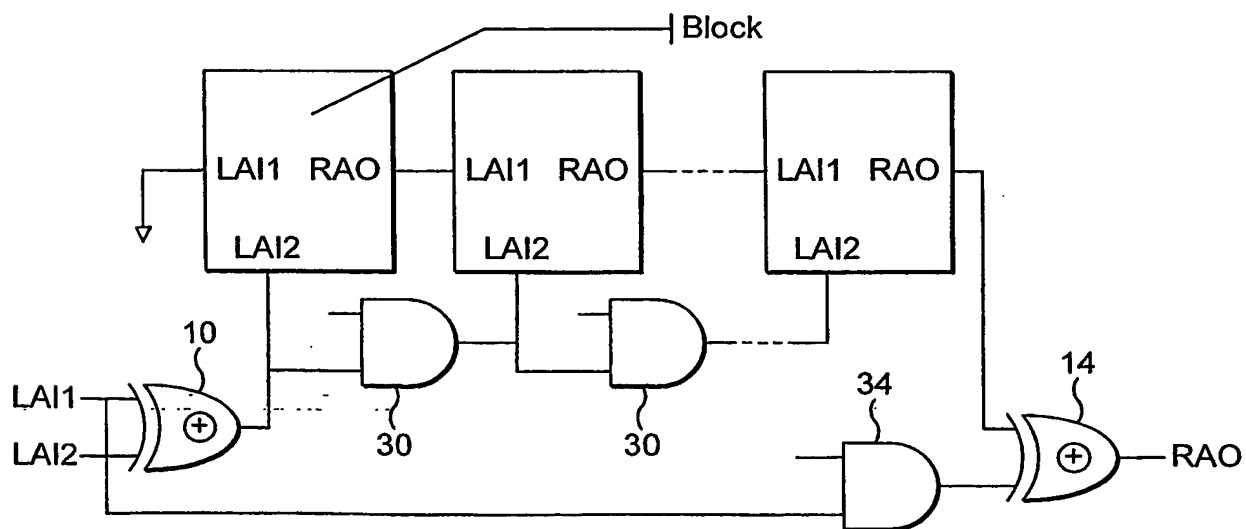


FIG. 7

4 / 8

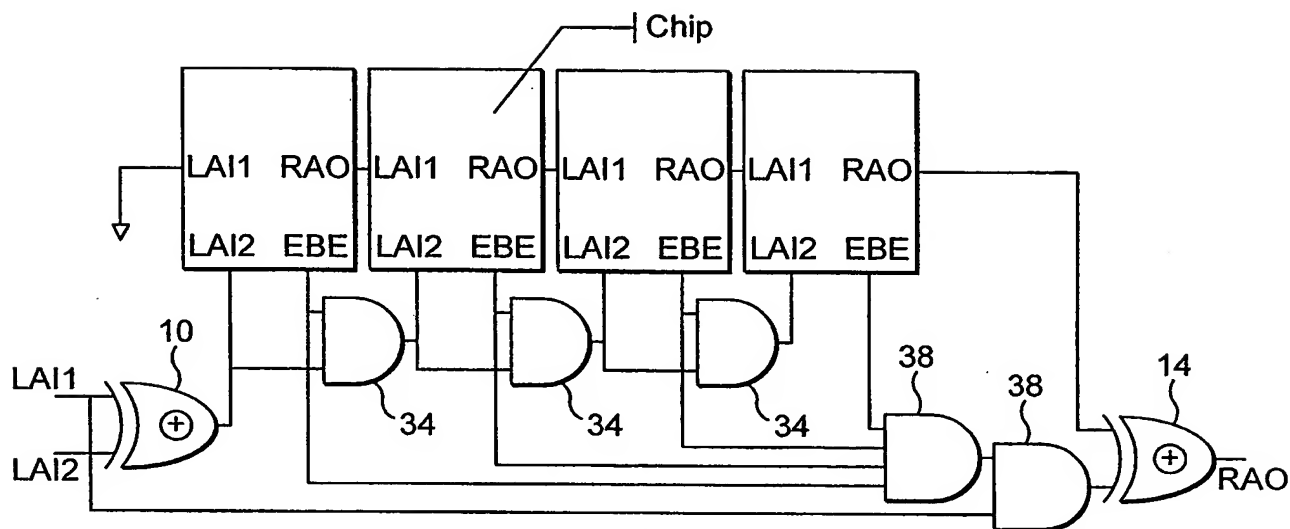


FIG. 8

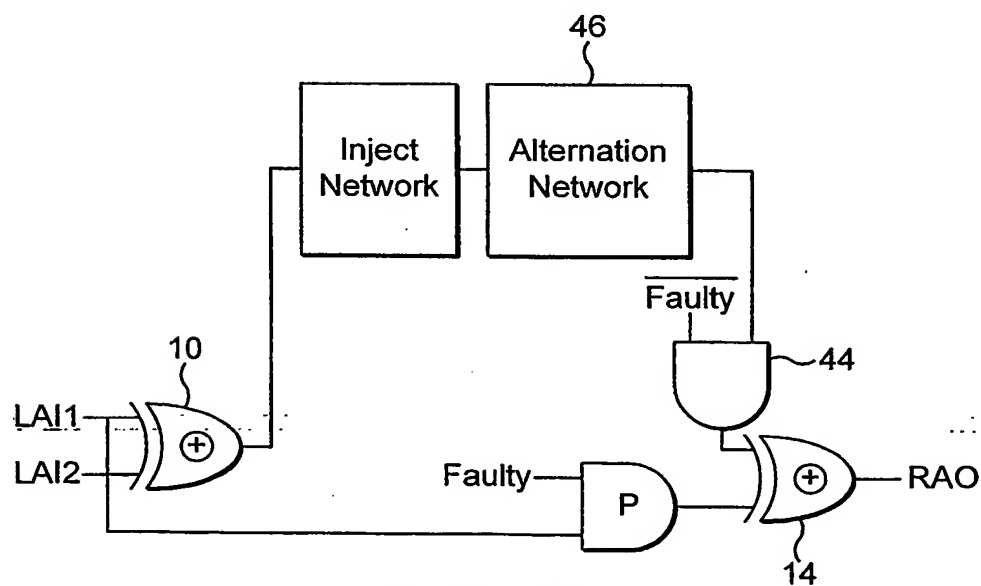


FIG. 10

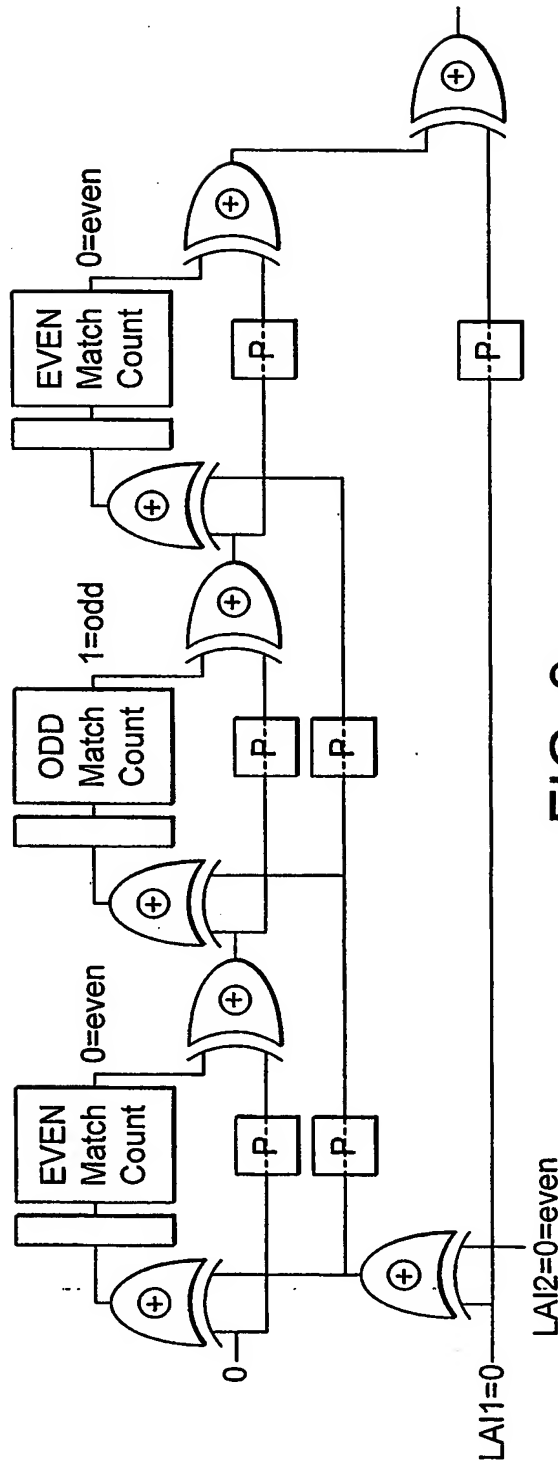


FIG. 9

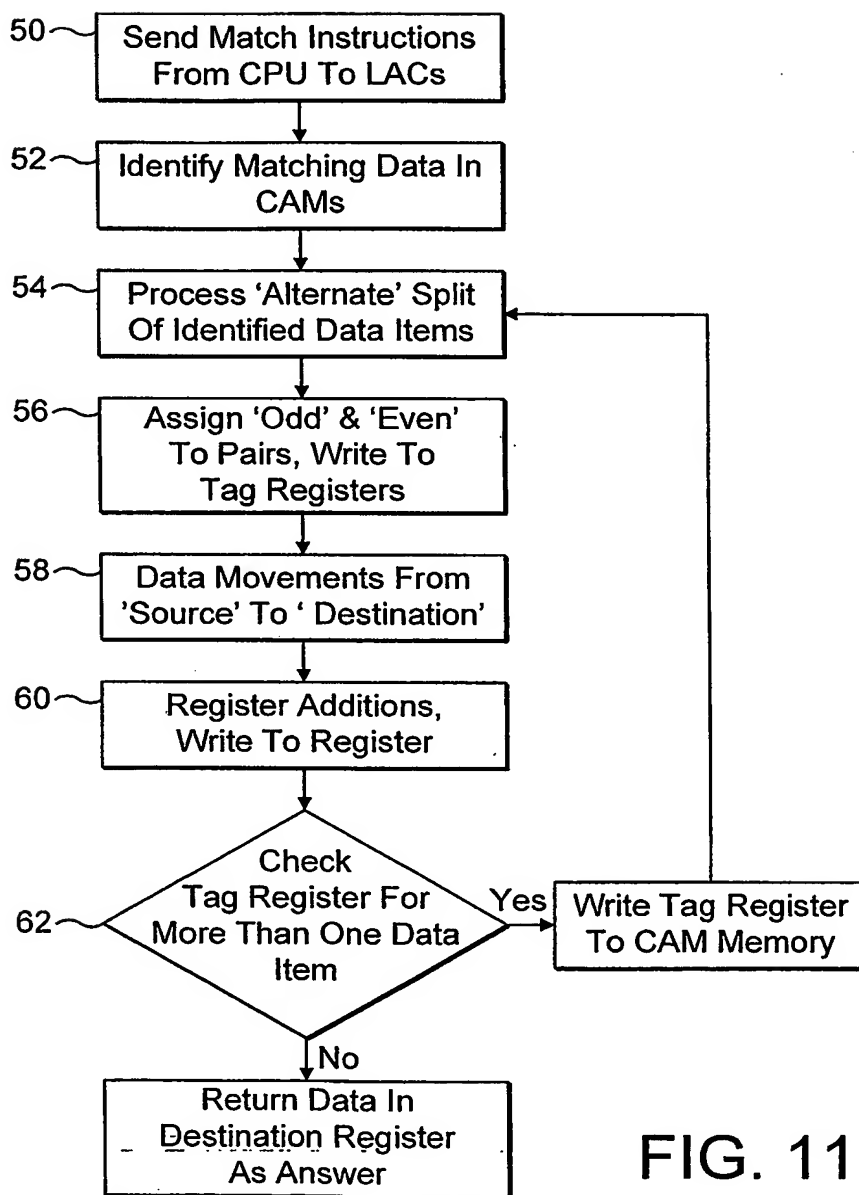


FIG. 11

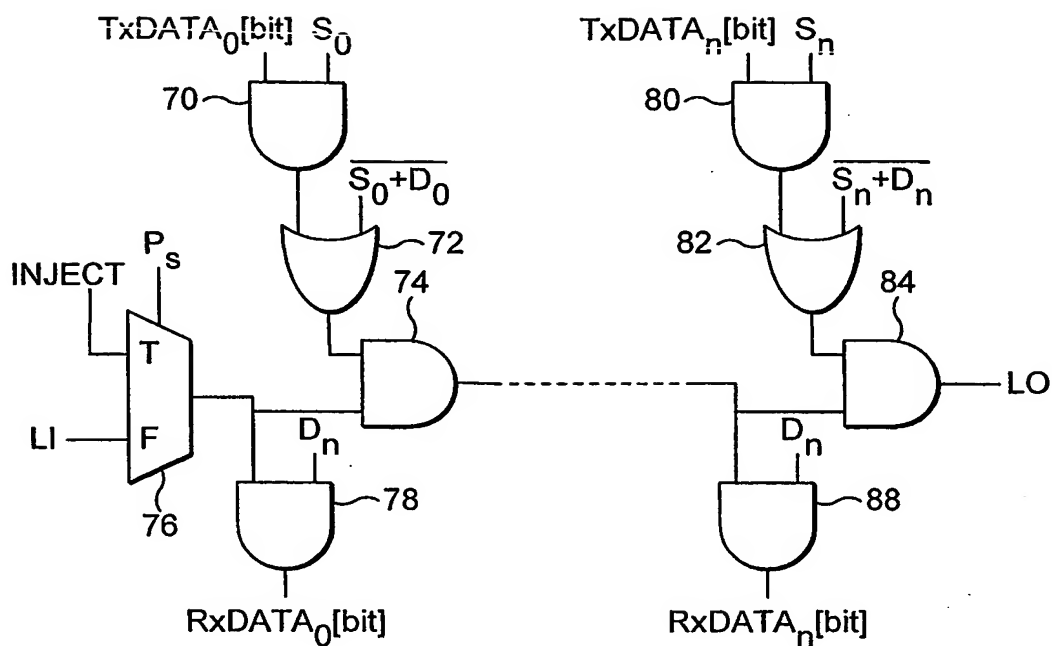


FIG. 12

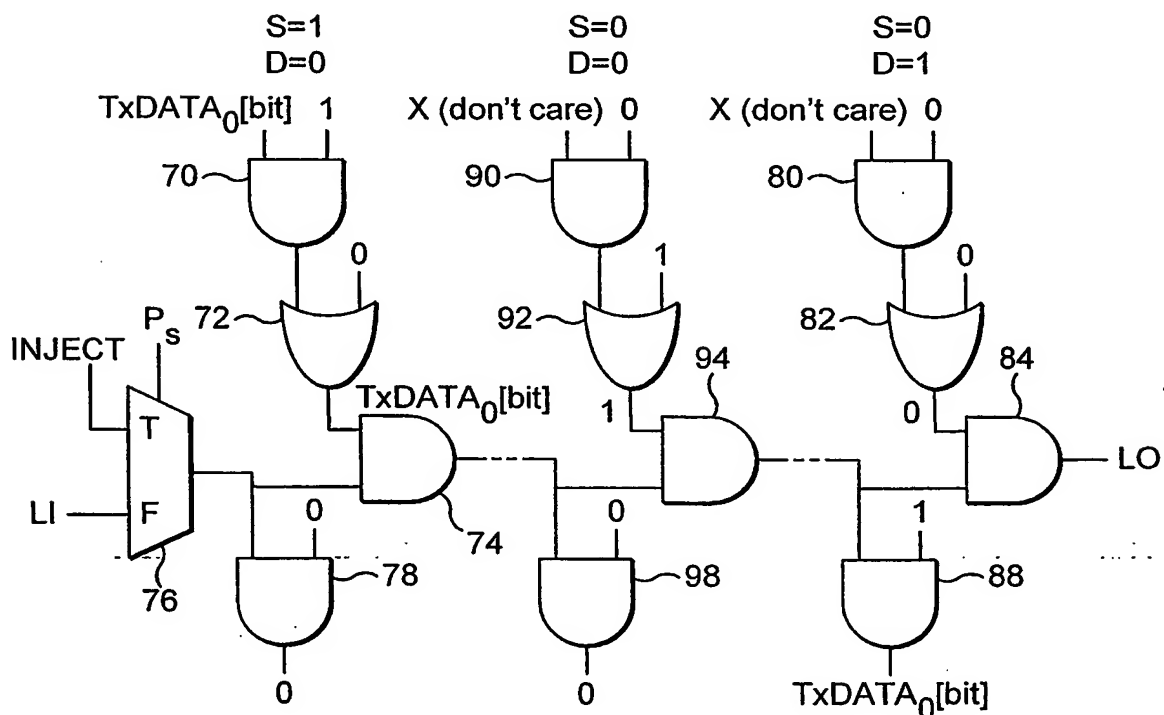
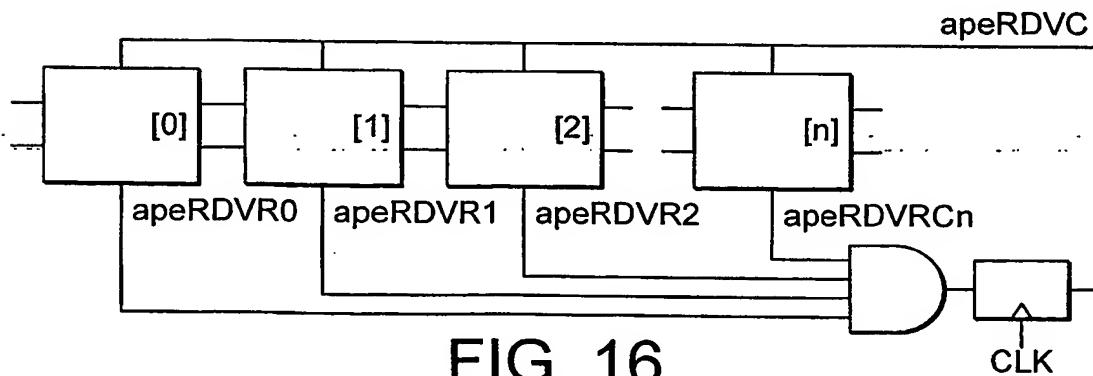
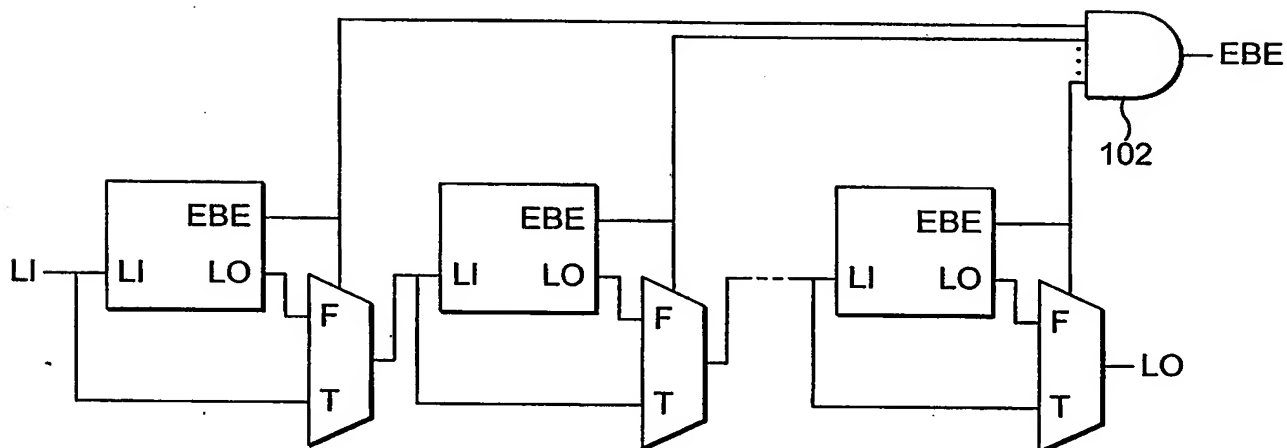
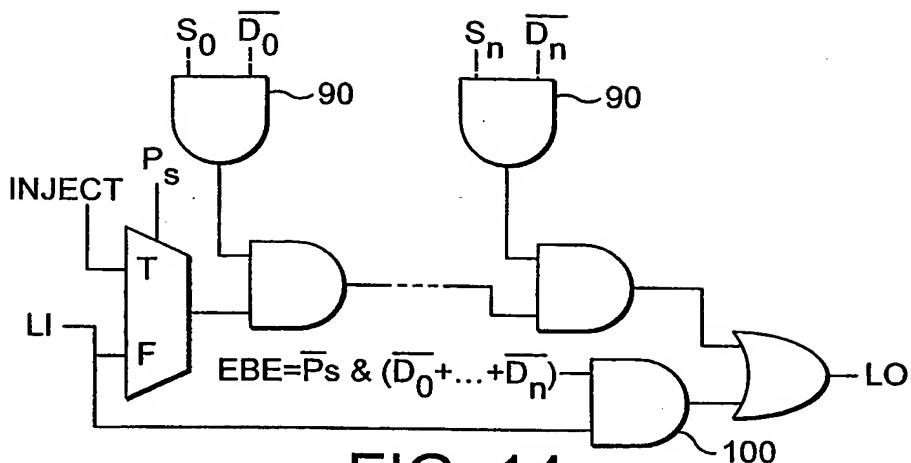


FIG. 13



INTERNATIONAL SEARCH REPORT

International Application No

PCT/GB 03/05532

A. CLASSIFICATION OF SUBJECT MATTER
IPC 7 G06F15/80

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, INSPEC, IBM-TDB

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>MARRIOTT A P ET AL: "Towards the automated design of application specific array processors (ASAPs)" APPLICATION SPECIFIC ARRAY PROCESSORS, 1990. PROCEEDINGS OF THE INTERNATIONAL CONFERENCE ON PRINCETON, NJ, USA 5-7 SEPT. 1990, LOS ALAMITOS, CA, USA, IEEE COMPUT. SOC, US, 5 September 1990 (1990-09-05), pages 414-425, XP010033964 ISBN: 0-8186-9089-5 page 415, line 1 - line 16 page 419, line 8 - page 420, line 16</p> <p style="text-align: center;">-/-</p>	1-37

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

* Special categories of cited documents :

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the international filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the international filing date but later than the priority date claimed

- *T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
- *Z* document member of the same patent family

Date of the actual completion of the international search

30 March 2004

Date of mailing of the international search report

08/04/2004

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Thibaudeau, J

INTERNATIONAL SEARCH REPORT

International Application No

PCT/GB 03/05532

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	GROSSPIETSCH K E ET AL: "THE ASSOCIATIVE PROCESSOR SYSTEM CAPRA: ARCHITECTURE AND APPLICATIONS" IEEE MICRO, IEEE INC. NEW YORK, US, vol. 12, no. 6, 1 December 1992 (1992-12-01), pages 58-67, XP000334759 ISSN: 0272-1732 page 64, left-hand column, line 2 -right-hand column, line 7; figure 5	1-37
A	US 5 974 521 A (AKERIB AVIDAN) 26 October 1999 (1999-10-26) column 20, line 1 -column 21, line 3	1-37